

# Logistic regression as a Neural Network

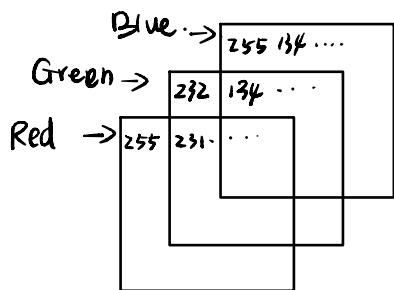
- Problem setting =

logistic regression is a type of binary classification.

Sample problem: classify image as 1 (Cat) or 2 (non-Cat)

Image size:  $64 \times 64$ , with 3 channels.

We usually layout each single input training data as a flattened 1-D array,



$$\text{In this case: } x = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 232 \\ 134 \\ \cdots \\ 255 \\ 231 \\ \cdots \end{bmatrix} \quad n = nx = 64 \times 64 \times 3 = 12288.$$

Our goal (model):  $x \rightarrow y$ ; Given any  $x$  ( $n_x = 12288$ ), predict  $y$  (1 or 0, stands for Cat or non-Cat).

- Notation.

$m$  training examples:  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

A single training example:  $(x, y)$ ,  $x \in \mathbb{R}^{n_x}$ ,  $y \in \{0, 1\}$

$M = M_{\text{train}} = \# \text{ of training data}$ ,  $M_{\text{test}} = \# \text{ test example}$ .

2 ways to represent training data / test data:

$$\textcircled{1} \quad X = \left[ \begin{array}{c|c|c|c} & & & \\ \hline x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ \hline & & & \end{array} \right] \quad \begin{matrix} \uparrow \\ n_x (\text{size of input data}) \end{matrix}$$

$$\textcircled{2} \quad \left[ \begin{array}{c} x^{(1)\top} \\ \hline x^{(2)\top} \\ \hline \vdots \\ \hline x^{(m)\top} \end{array} \right]$$

In order to reach the full power of neural network, we would normally use the first representation for data (input, output, training, test). More specifically, all data will be treated as **column vectors**.

In the previous example:  $X \in \mathbb{R}^{n \times m}$

$$X.\text{shape} = (\underline{n_x}, m).$$

$$Y = [y^{(1)}, y^{(2)} \dots y^{(m)}]$$

$$Y \in \mathbb{R}^{1 \times m}$$

$$Y.\text{shape} = (1, m)$$

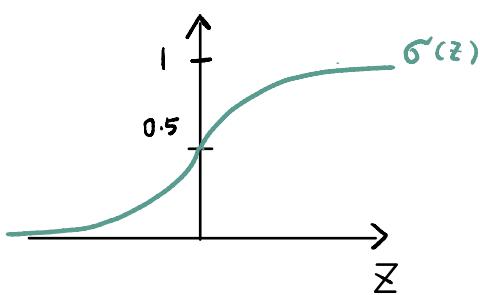
Formal Definition of Logistic Regression.

Given  $x$ , input, want  $\hat{y} = P(Y=1 | x) \in (0, 1)$ .

Parameters:  $w \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$

$$\Rightarrow \hat{y} = \frac{w^T x + b}{\text{we need } \hat{y} \text{ to be a number between 0 and 1.}}$$

$$\Rightarrow \hat{y} = \sigma(w^T x + b)$$



$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\text{if } z \text{ large, } \sigma(z) \approx \frac{1}{1+0} = 1$$

if  $z$  small (large negative)

$$\sigma(z) = \frac{1}{1+\infty} = 0$$

lost function and cost function:

problem:  $\hat{y} = \sigma(w^T x + b)$ , where  $\sigma(z) = \frac{1}{1+e^{-z}}$

Given  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$ , want  $\hat{y}^{(i)} \approx y^{(i)}$

We need a loss function to quantify how well does the model predict on the input data.

loss function: high loss, bad prediction.  
low loss, good prediction.

Any function of  $l(\hat{y}, y)$ , can be loss functions.  
However, some loss functions are not suitable for training neural networks.

e.g.  $\sum (y - \hat{y})^2$ : this is not a convex function,  
thus it has a lot of local mins,  
making it hard or impossible to work well with gradient descent.

Example of good loss function:

$$l(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

if  $y=1$ ,  $l(\hat{y}, y) = -\log \hat{y}$ ,  $\hat{y} \uparrow$ ,  $l(\hat{y}, y) \downarrow$  ✓

if  $y=0$ ,  $l(\hat{y}, y) = -\log(1-\hat{y})$ ,  $\hat{y} \downarrow$   $l(\hat{y}, y) \downarrow$  ✓

Cost function:

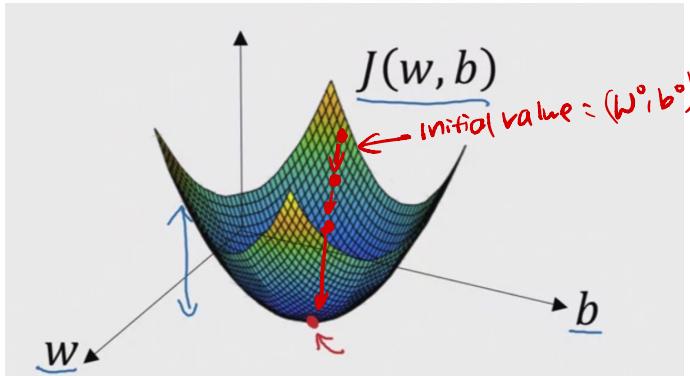
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}^{(i)}, y^{(i)})$$

Cost function gets cost of my parameters, in training, we wish to find the parameters  $w$  and  $b$  that minimize cost function.

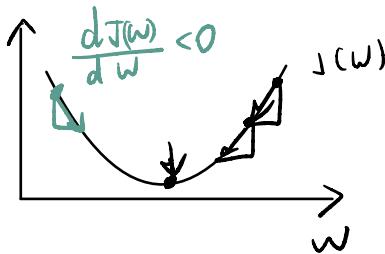
## Gradient Descent:

This is the technique that gradually adjusts training parameters and brings cost function to a local minimum.

Cost function for logistic regression looks like below.



Simpler case: ignoring  $b$  for now,



Repeat {  
 $w := w - \alpha$

$$\frac{dJ(w)}{dw}$$

Learning rate.  
 Gradient of cost  
 function at the  
 moment.

Back to logistic regression:

$$J(w, b)$$

$$w := w - \alpha$$

$$\frac{dJ(w, b)}{dw}$$

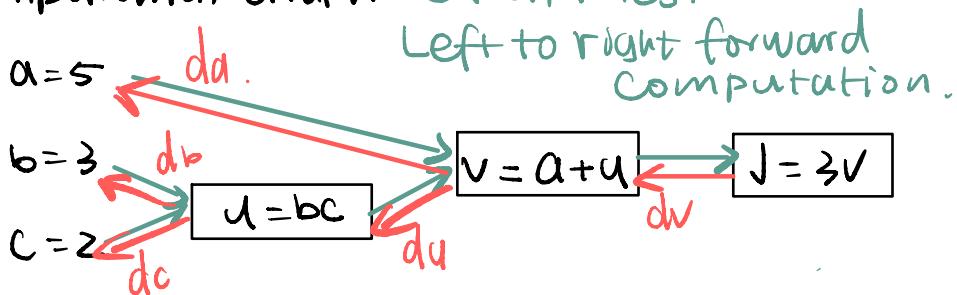
Sometimes  
 these are  
 written as  $\partial$ , for  
 partial derivative

$$b := b - \alpha \frac{dJ(w, b)}{db}$$

# Derivatives with a Computation Graph.

$$J(a, b, c) = 3(a+bc)$$

Computation Graph: Green lines.



Left to right forward computation.

Red lines: right to left, backward derivatives.

$da, db, dc, du, dv$  are abbs for partial derivatives

$$\frac{\partial J}{\partial *}$$

Multivariable calculus calculation:

$$dv = 3$$

$$du = \frac{\partial J}{\partial u} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} = 3 \cdot 1 = 3$$

$$da = \frac{\partial J}{\partial a} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial a} = 3 \cdot 1 = 3$$

$$db = \frac{\partial J}{\partial b} = \frac{\partial J}{\partial u} \cdot \frac{\partial u}{\partial b} = 3c$$

$$dc = \frac{\partial J}{\partial c} = \frac{\partial J}{\partial u} \cdot \frac{\partial u}{\partial c} = 3b$$

$da, db, dc$  are also gradients (directions at which  $a, b, c$  can change  $J$  the most).

logistic regression derivatives.

Assume each input has  $\geq$  features, i.e.  $x \in \mathbb{R}^2$ .

Diagram illustrating the forward pass and backpropagation for logistic regression:

Forward pass:  
Inputs:  $x_1, x_2, w_1, w_2, b$  feed into  $z = w_1 x_1 + w_2 x_2 + b$ .  
 $g(z) = 1 / (1 + e^{-z})$  is the sigmoid function.  
 $a = g(z)$  is the predicted probability.  
 $L(a, y)$  is the loss function.

Backpropagation (derivative calculations):

$$dz = \frac{dL(a, y)}{d(z)}$$
$$= -(y \log a + (1-y) \log(1-a))$$
$$dw_1 = x_1 dz = \frac{dL(a, y)}{da} \cdot \frac{da}{dz}$$
$$dw_2 = x_2 dz$$
$$db = dz$$
$$\frac{da}{da} = \frac{dL(a, y)}{da}$$
$$= -\frac{y}{a} + \frac{1-y}{1-a}$$

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

logistic regression on  $m$  examples

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \underbrace{\sum_{i=1}^m \frac{\partial}{\partial w_i} L(a^{(i)}, y^{(i)})}_{\text{d}w_1^{(i)} = (a^{(i)}, y^{(i)})}$$

Formal Algorithm.

$$J = 0; dw_1 = 0; dw_2 = 0; db = 0$$

for  $i = 1$  to  $m$ :

$$z^{(i)} = w^T x^{(i)} + b \quad \text{loss calculation}$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)} \quad \text{Derivative Calculation.}$$

$$dw_1 += x_1^{(i)} dz^{(i)}; dw_2 += x_2^{(i)} dz^{(i)}; db += dz^{(i)}$$

$$\downarrow i = m$$

$$dw_1 / = m; dw_2 / = m; db / = m$$

We have too many for loops which is not very efficient. Thus, we need Vectorization.